

## Aberystwyth University

### *Homology Induction: the use of machine learning to improve sequence similarity searches*

Karwath, Andreas; King, Ross Donald

*Published in:*  
BMC Bioinformatics

*DOI:*  
[10.1186/1471-2105-3-11](https://doi.org/10.1186/1471-2105-3-11)

*Publication date:*  
2002

*Citation for published version (APA):*  
Karwath, A., & King, R. D. (2002). Homology Induction: the use of machine learning to improve sequence similarity searches. *BMC Bioinformatics*, 3(11). <https://doi.org/10.1186/1471-2105-3-11>

#### **General rights**

Copyright and moral rights for the publications made accessible in the Aberystwyth Research Portal (the Institutional Repository) are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Aberystwyth Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Aberystwyth Research Portal

#### **Take down policy**

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

tel: +44 1970 62 2400  
email: [is@aber.ac.uk](mailto:is@aber.ac.uk)

species(Name, Species).

Describes the organisms species declaration in one string, e.g.

species(p02102, capra\_hircus\_goat).

classification(Name, Classification\_class).

This declaration defines the peptides' organisms origin classification. The classification is divided from top to bottom of the phylogenic tree, starting with the most general classification and ending with the most specific classification. To improve computational efficiency, each leaf of the tree is an entry in the database. The following example classifications are taken from the SWISS-PROT protein P02102.

- classification(p02102, eukaryota).
- classification(p02102, metazoa).
- classification(p02102, chordata).
- classification(p02102, vertebrata).
- classification(p02102, tetrapoda).
- classification(p02102, mammalia).
- classification(p02102, eutheria).
- classification(p02102, artiodactyla).

desc(Name, Description).

Defines the protein's description word by word.

keyword(Name, Keyword).

Defines the proteins keyword.

mol\_wt\_rule(Name, Mol\_weight\_interval).

The relative molecular weight of the protein.

db\_ref(Name, Database\_idenfier, Primary\_idenfier, Secondary\_idenfier, Status).

Defines a database reference to a specific database and its entry.

domain\_rule(Name, Query\_start\_interval, Query\_end\_interval, Target\_start\_interval, Target\_end\_interval).

Defines the relative hit taken from the local PSI-BLAST alignment.

amino\_acid\_ratio\_rule(Name, Amino\_acid, Percentage\_interval).

Defines the content of a particular amino acid on a 1-10 scale.

amino\_acid\_pair\_ratio\_rule(Name, Amino\_acid, Amino\_acid, Promille\_interval).

Measures the pairwise amino acid content on a 1-10 scale.

seq\_length\_rule(Name, Seq\_length\_interval).

Measures the total number of amino acids in the sequence on a scale from 1 to 10.

signalip1\_rule(Name, Signalip1\_interval).

Defines the predicted cleavage sites found using the SignalIP program (Nielsen et. al, 1997) server. Signalip1 is the protein's maximum cleavage site score.

signalip2\_rule(Name, Signalip2\_interval).

Signalip2 is the protein's maximum combined cleavage site score

signalip3\_rule(Name, Signalip3\_interval).

Signalip3 is the protein's maximum signal peptide score

signalip4\_rule(Name, Signalip4\_interval1, Signalip4\_interval2).

The signalip4 interval is the proteins most likely cleavage site

hydro\_cons\_rule(Name, hc\_interval1, hc\_interval2, hc\_interval3, hc\_interval4).

hc\_interval1 defines the average hydrophobic moment assuming  $\alpha$ -helix,  
 hc\_interval2 the average hydrophobic moment assuming  $\beta$ -strand, and  
 hc\_interval3 the average hydropathy using the Kyte-Doolittle hydrophilicity  
 scale  
 sec\_struc\_rule(Name, Sec\_position\_interval, Secondary\_structure,  
 Sec\_length\_interval).  
 Defines the position, length, and type of predicted secondary structure.  
 sec\_struc\_alpha\_rule(Name, Sec\_position\_alpha\_interval, Sec\_length\_alpha\_interval).  
 Defines the position and length of predicted secondary structure of type  $\alpha$ .  
 sec\_struc\_beta\_rule(Name, Sec\_position\_beta\_interval, Sec\_length\_beta\_interval).  
 Defines the position and length of predicted secondary structure of type  $\beta$ .  
 sec\_struc\_coil\_rule(Name, sec\_position\_coil\_interval, sec\_length\_coil\_interval).  
 Defines the position and length of predicted secondary structure of type coil.  
 sec\_struc\_distribution\_rule(Name, Secondary\_structure, Sec\_dist\_interval).  
 Measures the distribution of all three types of secondary structure predictions  
 sec\_struc\_conf\_rule(Name, Sec\_conf\_interval).  
 Measures the confidence in the over all secondary structure prediction.  
 sec\_struc\_conf\_alpha\_rule(Name, Sec\_conf\_alpha\_interval).  
 Measures the confidence in the  $\alpha$  secondary structure predictions.  
 sec\_struc\_conf\_beta\_rule(Name, Sec\_conf\_beta\_interval).  
 Measures the confidence in the  $\beta$  secondary structure predictions.  
 sec\_struc\_conf\_coil\_rule(Name, Sec\_conf\_coil\_interval).  
 Measures the confidence in the coil secondary structure predictions.

Table 1. The definition of types of data held in the logical database. The  
 data was obtained from a wide variety of bioinformatic sources. This  
 information was selected for relevance to the detection of homology. For  
 each protein we collected information for each type of data if it was available.